

	age	duration	campaign	pdays	previous	emp.var.rate	cons.price.idx	cons.conf.idx	euribor3m
count	37803.000000	37803.000000	37803.000000	37803.000000	37803.000000	37803.000000	37803.000000	37803.000000	37803.000000
mean	40.007142	263.432217	1.941751	960.176653	0.181785	0.028646	93.560179	-40.485386	1.56913
std	10.458787	259.252134	1.125077	982.472659	0.506630	1.583386	0.581213	4.693541	3.74831
min	17.000000	0.000000	1.000000	0.000000	0.000000	-3.400000	92.201000	-50.800000	0.6340X
25%	32.000000	107.000000	1.000000	999.000000	0.000000	-1.800000	93.075000	-42.700000	1.3340X
50%	38.000000	185.000000	2.000000	999.000000	0.000000	1.000000	93.444000	-41.800000	4.8570X
75%	47.000000	324.000000	3.000000	999.000000	0.000000	1.400000	93.994000	-36.400000	4.9610X
max	90.000000	4918.000000	5.000000	999.000000	7.000000	1.400000	94.767000	-26.900000	5.0450X

```
In [715]: encoder_df = encoder_df.drop('duration', axis=1)
```

```
In [717]: encoder_df.reset_index(inplace=True)
```

	index	age	job	marital	education	default	housing	loan	contact	month	...	campaign	pdays	previous	outcome	emp.var.rate
0	0	56	3	1	0	0	0	0	1	6	...	1	999	0	1	
1	1	57	7	1	3	1	0	0	1	6	...	1	999	0	1	
2	2	37	7	1	3	0	2	0	1	6	...	1	999	0	1	
3	3	40	0	1	1	0	0	0	1	6	...	1	999	0	1	
4	4	56	7	1	3	0	0	2	1	6	...	1	999	0	1	
...	
37798	41183	73	5	1	5	0	2	0	0	7	...	1	999	0	1	-
37799	41184	48	1	1	5	0	0	0	0	7	...	1	999	0	1	-
37800	41185	56	5	1	6	0	2	0	0	7	...	2	999	0	1	-
37801	41186	44	9	1	5	0	0	0	0	7	...	1	999	0	1	-
37802	41187	74	5	1	5	0	2	0	0	7	...	3	999	1	0	-

37803 rows × 21 columns

```
In [718]: #standardizing the dataset
scaler = preprocessing.StandardScaler()
std_df = scaler.fit_transform(encoder_df.drop(['job', 'marital', 'education', 'default', 'housing', 'loan', 'co
scaled_df = pd.DataFrame(std_df, columns=['age', 'pdays', 'cons.price.idx', 'cons.conf.idx', 'euribor3m', 'nr.empl
```

```
In [719]: scaled_df = pd.DataFrame(std_df, columns=['age', 'pdays', 'cons.price.idx', 'cons.conf.idx', 'euribor3m', 'nr.empl
```

```
In [720]: scaled_df[['job']] = encoder_df[['job']].values
scaled_df[['marital']] = encoder_df[['marital']].values
scaled_df[['education']] = encoder_df[['education']].values
scaled_df[['default']] = encoder_df[['default']].values
scaled_df[['housing']] = encoder_df[['housing']].values
scaled_df[['loan']] = encoder_df[['loan']].values
scaled_df[['contact']] = encoder_df[['contact']].values
scaled_df[['month']] = encoder_df[['month']].values
scaled_df[['day.of.week']] = encoder_df[['day.of.week']].values
scaled_df[['campaign']] = encoder_df[['campaign']].values
scaled_df[['previous']] = encoder_df[['previous']].values
scaled_df[['outcome']] = encoder_df[['outcome']].values
scaled_df[['emp.var.rate']] = encoder_df[['emp.var.rate']].values
scaled_df[['y']] = encoder_df[['y']].values
```

```
In [721]: scaled_df.head()
```

	age	pdays	cons.price.idx	cons.conf.idx	euribor3m	nr.employed	job	marital	education	default	housing	loan	contact	mo
0	1.529151	0.207172	0.746417	0.870441	0.736639	0.360955	3	1	0	0	0	0	1	
1	1.624766	0.207172	0.746417	0.870441	0.736639	0.360955	7	1	3	1	0	0	1	
2	-0.287527	0.207172	0.746417	0.870441	0.736639	0.360955	7	1	3	0	2	0	1	
3	-0.005683	0.207172	0.746417	0.870441	0.736639	0.360955	0	1	1	0	0	0	1	
4	1.529151	0.207172	0.746417	0.870441	0.736639	0.360955	7	1	3	0	0	2	1	

```
In [722]: scaled_df.describe()
```

	age	pdays	cons.price.idx	cons.conf.idx	euribor3m	nr.employed	job	marital	education	default	housing	loan	contact	mo
count	3.780300e+04	3.780300e+04	3.780300e+04	3.780300e+04	3.780300e+04	3.780300e+04	3.780300e+04	3.780300e+04	3.780300e+04	3.780300e+04	3.780300e+04	3.780300e+04	3.780300e+04	3.780300e+04
mean	3.247938e+16	-7.217639e-17	-1.814033e-14	0.000000	2.405880e-17	4.811759e-16	1.941751	0.181785	0.028646	0.506630	1.583386	0.581213	4.693541	3.74831
std	1.000013e+00	1.000013e+00	1.000013e+00	1.000013	1.000013e+00	1.000013e+00	1.125077	0.506630	1.583386	0.000000	-3.400000	-50.800000	-42.700000	1.3340X
25%	-7.956001e-01	2.071171e-01	-8.347802e-01	-0.471847	-1.678834e+00	-2.752819e-01	1.000000	0.000000	-1.800000	0.000000	-1.800000	0.000000	-1.800000	-1.800000
50%	-1.919122e-01	2.071171e-01	-1.998926e-01	-0.287009	7.366394e-01	3.609552e-01	2.000000	0.000000	1.000000	0.000000	1.000000	0.000000	1.000000	1.000000
75%	6.686196e-01	2.071171e-01	7.464167e-01	0.870441	7.961253e-01	6.686932e-01	3.000000	0.000000	1.400000	0.000000	1.400000	0.000000	1.400000	1.400000
max	5.544967e+00	2.071171e-01	2.076411e+00	2.894526	8.441717e-01	8.689632e-01	5.000000	7.000000	1.400000	1.400000	1.400000	1.400000	1.400000	1.400000

Problem 6: Train/Test Split

With your data prepared, split it into a train and test set.

```
In [723]: X = scaled_df.drop('y', axis=1)
y = scaled_df['y'].astype('int')
train(X.shape)
print(X.shape)
print(y.shape)
```

```
(37803, 19)
```

```
(37803, 1)
```

```
In [724]: X_train,X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=42)
```

```
In [725]: print(X_train.shape)
print(y_train.shape)
print(X_test.shape)
print(y_test.shape)
```

```
(25328, 19)
```

```
(12475, 1)
```

```
(12475, 19)
```

```
(12475, 1)
```

Problem 7: A Baseline Model

Before we build our first model, we want to establish a baseline. What is the baseline performance that our classifier should aim to beat?

```
In [726]: models.append(['Baseline'])
train_time = []
train_accuracy = []
test_accuracy = []
accuracy_score = []
AUC_score = []
```

```
In [727]: # Baseline Model
dummy_clf = DummyClassifier( random_state=42)
start_time = time()
dummy_clf.fit(X_train, y_train)
train_time.append(time() - start_time)
# Training
y_preds = dummy_clf.predict(X_train)
training_score = dummy_clf.score(X_train, y_train)
train_accuracy.append('training_score')
# Test score
y_preds_test = dummy_clf.predict(X_test)
test_score = dummy_clf.score(X_test, y_test)
test_accuracy.append(test_score)
accuracy_score.append('A/A')
AUC_score.append('N/A')
```

Problem 8: A Simple Model

Use Logistic Regression to build a basic model on your data.

```
In [728]: # Logistics Regression Model
models.append('Logistics Regression')
log_reg = LogisticRegression(solver='liblinear', random_state=42)
start_time = time()
log_reg.fit(X_train, y_train)
train_time.append(time() - start_time)
y_preds_log_reg = log_reg.predict(X_test)
```

Problem 9: Score the Model

What is the accuracy of your model?

```
In [729]: train_score = log_reg.score(X_train, y_train)
train_accuracy.append(train_score)
test_accuracy.append(log_reg.score(X_test, y_test))
#accuracy scores
accuracy_score.append(metrics.accuracy_score(y_test, y_preds_log_reg))
fpr, tpr, _thresholds = metrics.roc_curve(y_test, y_preds_log_reg)
AUC_score.append(metrics.auc(fpr, tpr))
```

```
In [730]: print('Precision : %.3f'%precision_score(y_test, y_preds_log_reg))
print('Recall : %.3f'%recall_score(y_test, y_preds_log_reg))
print('F1 Score : %.3f'%f1_score(y_test, y_preds_log_reg))
```

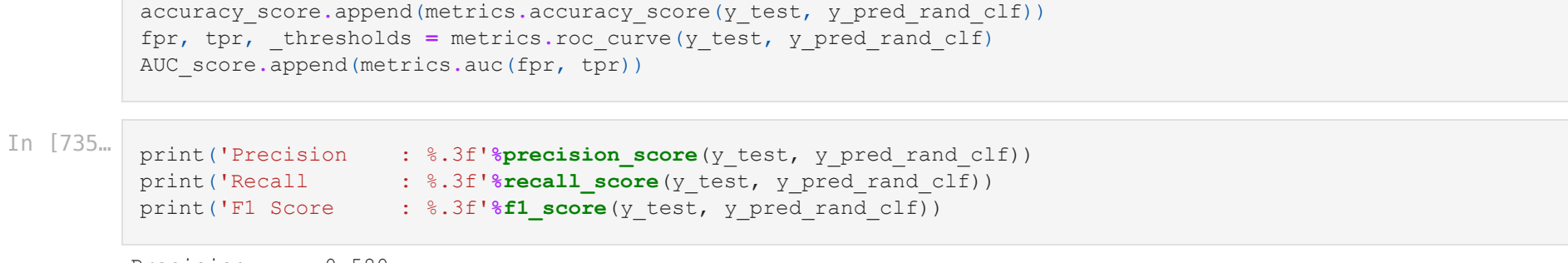
```
Precision : 0.665
Recall : 0.237
F1 Score : 0.349
```

```
In [731]: from sklearn.metrics import confusion_matrix
confusion_matrix = confusion_matrix(y_test, y_test_preds)
print(confusion_matrix)
```

```
[[10003 179]
 [1140 354]]
```

```
In [732]: disp = ConfusionMatrixDisplay(confusion_matrix=confusion_matrix)
disp.plot()
```

```
Out[732]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x146d51f0>
```



Problem 10: Model Comparisons

Now, we aim to compare the performance of the Logistic Regression model to our KNN algorithm, Decision Tree, and SVM models. Using the default settings for each of the models, fit and score each. Also, be sure to compare the fit time of each of the models. Present your findings in a DataFrame similar to that below:

	Model	Train Time	Train Accuracy	Test Accuracy
0	Baseline	0.007880	training_score	0.680240
1	Logistics Regression	0.164353	0.89711	0.894349
2	RandomForestClassifier	1.631834	0.994709	0.891062
3	KNeighborsClassifier	0.026783	0.909428	0.885291
4	SVC	19.737240	0.896162	0.894669

RandomForestClassifier

```
In [733]: models.append('RandomForestClassifier')
rand_clf = RandomForestClassifier(random_state=42)
start_time = time()
rand_clf.fit(X_train, y_train)
train_time.append(time() - start_time)
y_pred_rand_clf = rand_clf.predict(X_test)
```

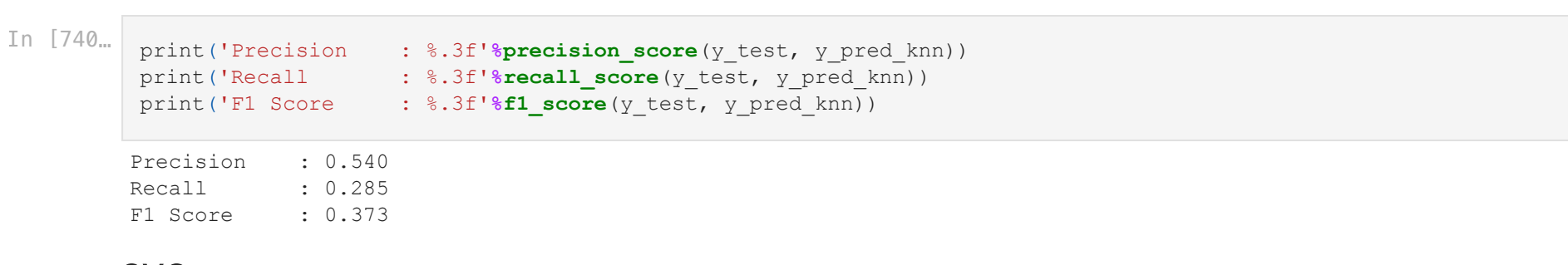
```
In [734]: # Scoring the model
train_score = rand_clf.score(X_train, y_train)
train_accuracy.append(train_score)
test_accuracy.append(rand_clf.score(X_test, y_test))
#accuracy scores
accuracy_score.append(metrics.accuracy_score(y_test, y_pred_rand_clf))
fpr, tpr, _thresholds = metrics.roc_curve(y_test, y_pred_rand_clf)
AUC_score.append(metrics.auc(fpr, tpr))
```

```
In [735]: print('Precision : %.3f'%precision_score(y_test, y_pred_rand_clf))
print('Recall : %.3f'%recall_score(y_test, y_pred_rand_clf))
print('F1 Score : %.3f'%f1_score(y_test, y_pred_rand_clf))
```

```
Precision : 0.580
Recall : 0.327
F1 Score : 0.418
```

```
In [736]: cm = metrics.confusion_matrix(y_test, y_pred_rand_clf)
disp = ConfusionMatrixDisplay(confusion_matrix=cm)
disp.plot()
```

```
Out[736]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x146d4dc0d>
```



```
In [737]: print('Precision : %.3f'%precision_score(y_test, y_pred_knn))
print('Recall : %.3f'%recall_score(y_test, y_pred_knn))
print('F1 Score : %.3f'%f1_score(y_test, y_pred_knn))
```

```
Precision : 0.540
Recall : 0.285
F1 Score : 0.373
```

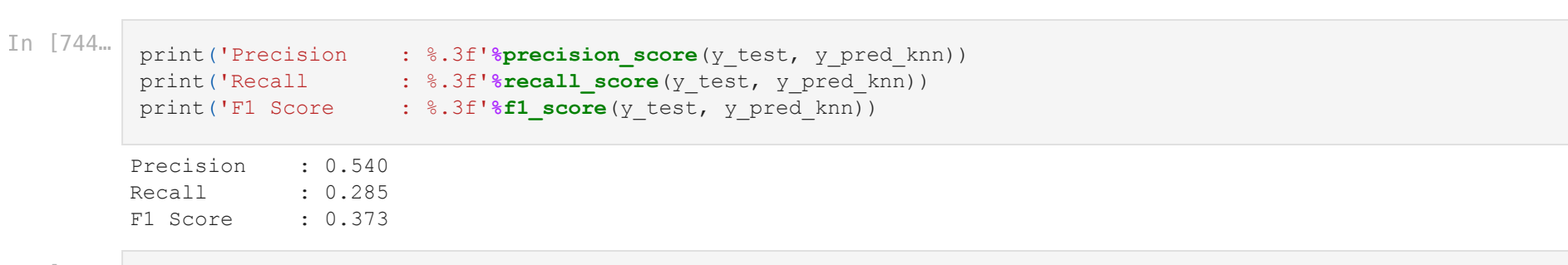
SVC

```
In [741]: models.append('SVC')
svc = SVC()
start_time = time()
svc.fit(X_train, y_train)
train_time.append(time() - start_time)
y_pred_svc = svc.predict(X_test)
```

```
In [742]: # Scoring the model
train_score = svc.score(X_train, y_train)
train_accuracy.append(train_score)
test_accuracy.append(svc.score(X_test, y_test))
#accuracy scores
accuracy_score.append(metrics.accuracy_score(y_test, y_pred_svc))
fpr, tpr, _thresholds = metrics.roc_curve(y_test, y_pred_svc)
AUC_score.append(metrics.auc(fpr, tpr))
```

```
In [743]: cm = metrics.confusion_matrix(y_test, y_pred_svc)
disp = ConfusionMatrixDisplay(confusion_matrix=cm)
disp.plot()
```

```
Out[743]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x146891460>
```



```
In [744]: print('Precision : %.3f'%precision_score(y_test, y_pred_knn))
print('Recall : %.3f'%recall_score(y_test, y_pred_knn))
print('F1 Score : %.3f'%f1_score(y_test, y_pred_knn))
```

```
Precision : 0.540
Recall : 0.285
F1 Score : 0.373
```

```
In [745]: model_artifacts = pd.DataFrame({'Model':models,
'Training Time':train_time,
'Training Score': train_accuracy,
'Test Score':test_accuracy,
'Accuracy': accuracy_score,
'AUC': AUC_score})
```

```
Out[745]:
```

	Model	Training Time	Training Score	Test Score	Accuracy	AUC
0	Baseline	0.007880	training_score	0.680240	N/A	N/A
1	Logistics Regression	0.164353	0.89711	0.894349	0.894349	0.610369
2	RandomForestClassifier	1.631834	0.994709	0.891062	0.891062	0.647336
3	KNeighborsClassifier	0.026783	0.909428	0.885291	0.885291	0.628042
4	SVC	19.737240	0.896162	0.894669	0.894669	0.598553

Problem 11: Improving the Model

Now that we have some basic models on the board, we want to try to improve these. Below, we list a few things to explore in this pursuit.

- More feature engineering and exploration. For example, should we keep the gender feature? Why or why not?
- Hyperparameter tuning and grid search. All of our models have additional hyperparameters to tune and explore. For example the number of neighbors in KNN or the maximum depth of a Decision Tree.
- Adjust your performance metric

Logistics Regression models training, test, AUC scores better than other models. Optimizing Logistics Regression models to improve further.

```
In [746]: grid={ "C":np.logspace(-3,3,7), "penalty":["l1","l2"]}
log_reg = LogisticRegression()
lgr_cv=GridSearchCV(lgr, grid, cv=10)
lgr_cv.fit(X_train, y_train)
```

```
Out[746]: GridSearchCV(cv=10,
estimator=LogisticRegression(random_state=41, solver='liblinear'),
param_grid={ 'C': array([1.e+03, 1.e+02, 1.e+01, 1.e+00, 1.e+01, 1.e+02, 1.e+03]),
'penalty': ['l1', 'l2']})
```

```
In [747]: lgr_cv.best_params_
```

```
{'C': 0.1, 'penalty': 'l1'}
```

```
In [748]: #Adjust Performance Metrics
lgr2=LogisticRegression(C=0.1, penalty='l2', solver='liblinear')
start_time = time()
lgr2.fit(X_train,y_train)
train_time = time() - start_time
```

```
In [749]: y_pred = lgr2.predict(X_test)
```

```
In [750]: cm = metrics.confusion_matrix(y_test, y_pred)
```

```
In [751]: disp = ConfusionMatrixDisplay(confusion_matrix=cm)
disp.plot()
```

```
Out[751]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x15c6a6d0>
```



```
In [752]: print('Precision : %.3f'%precision_score(y_test, y_pred))
print('Recall : %.3f'%recall_score(y_test, y_pred))
print('F1 Score : %.3f'%f1_score(y_test, y_pred))
```

```
Precision : 0.637
Recall : 0.234
F1 Score : 0.345
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

Questions